

STAT 542: Statistical Learning

Penalized Linear Regression: Part I

Ruoqing Zhu, Ph.D. <rqzhu@illinois.edu>

Course Website: <https://teazrq.github.io/stat542/>

February 12, 2022

Department of Statistics
University of Illinois at Urbana-Champaign

Shrinkage Methods

- Best subset selection
 - Computationally expensive
 - Not feasible when p is large
- Forward/backward selection
 - No guarantee to find the best global submodel
 - The selection process is discrete (“add” or “drop”), often leads to high variance.
- Shrinkage methods
 - A continuous process, does not suffer from high variability

Motivation

- The OLS estimator is a linear function of \mathbf{y} , and it is the BLUE.
- But there can be (and often exist) biased estimators with smaller variance
- Recall that the **prediction accuracy** is

$$\text{Bias}^2 + \text{Variance} + \text{Irreducible Error}$$

and choosing estimators often involves the bias-variance trade-off.

- Generally, by **regularizing** (shrinking, penalizing) the estimator in some way, its variance can be reduced; if the corresponding increase in bias is small, we have better prediction accuracy

- Part I
 - l_2 penalty: Ridge regression
 - l_1 penalty: Lasso
 - Connecting the two: Elastic net; Bridge penalty
- Part II
 - Bias reduction: adaptive Lasso, SCAD, MCP
 - Consistency of penalized methods
 - Penalties for special data structures: grouped lasso, fused lasso

Ridge Regression

A Motivating Example

```
1 > library(MASS)
2 > set.seed(1)
3 > n = 30
4 >
5 > # highly correlated variables
6 > X = mvrnorm(n, c(0, 0), matrix(c(1,0.999, 0.999, 1), 2,2))
7 > y = rnorm(n, mean=1 + X[,1] + X[,2])
8 >
9 > # compare parameter estimates
10 > summary(lm(y~X))$coef
11      Estimate Std. Error  t value    Pr(>|t|)
12 (Intercept)  1.038007  0.1647551  6.300302 9.627026e-07
13 X1          -11.272638  4.6402098 -2.429338 2.205727e-02
14 X2           13.265586  4.6315269  2.864193 7.993486e-03
15 > lm.ridge(y~X, lambda=5)
16      X1      X2
17 1.1214448 0.8770568 0.9836474
```

Penalizing the square of the coefficients

$$\begin{aligned}\hat{\beta}^{\text{ridge}} &= \arg \min_{\beta} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|^2 & (1) \\ &= \arg \min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|^2 + n\lambda \|\beta\|^2\end{aligned}$$

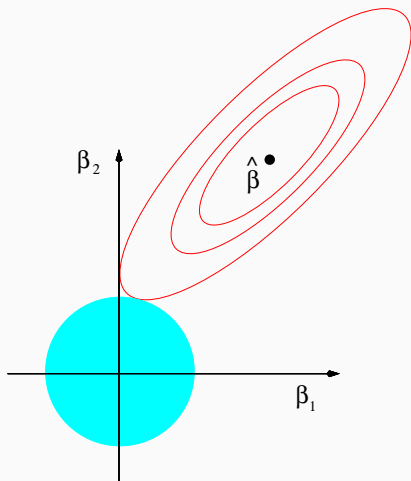
- Hoerl and Kennard (1970); Tikhonov (1943)
- $\lambda \geq 0$ is a **tuning parameter** (penalty level), it controls the amount of shrinkage.
- The coefficients $\hat{\beta}^{\text{ridge}}$ are shrunken towards 0.

An **equivalent formulation** is given by

$$\begin{aligned} & \underset{\beta}{\text{minimize}} && \sum_{i=1}^n \left(y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \\ & \text{subject to} && \sum_{j=1}^p \beta_j^2 \leq s \end{aligned}$$

- There is a **one-to-one correspondence** between the parameters λ and s
- This is due to the KKT conditions

Ridge Regression



Ridge constrained solution

Ridge Regression

- Ridge regression is mainly used to address multi-collinearity problem in high-dimensional data
- When there are many correlated variables, a wildly large positive coefficient on one variable can be canceled by a similarly large negative coefficient on its correlated cousin.
- Ridge regression alleviate this problem by imposing a size constraint

Ridge Regression

- How to derive the solution $\hat{\beta}^{\text{ridge}}$
- Degrees of freedom
- Tuning parameter selection
- Connections with other methods

Solution for Ridge Regression

- For a fixed tuning parameter λ , we want to minimize

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + n\lambda\boldsymbol{\beta}^\top\boldsymbol{\beta}$$

- Take derivative with respect to $\boldsymbol{\beta}$ and set to zero, we have the solution of the Ridge regression

$$\hat{\boldsymbol{\beta}}^{\text{ridge}} = (\mathbf{X}^\top\mathbf{X} + n\lambda\mathbf{I})^{-1}\mathbf{X}^\top\mathbf{y}$$

- $\hat{\boldsymbol{\beta}}^{\text{ridge}}$ is still a **linear** estimator

Solution for Ridge Regression

- This is similar to the ordinary least squares solution, but with the addition of a “ridge” down the diagonal
- $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$ is **always invertible**, hence $\hat{\beta}^{\text{ridge}}$ is unique
- As $\lambda \rightarrow 0$, $\hat{\beta}^{\text{ridge}} \rightarrow \hat{\beta}^{\text{ols}}$
- As $\lambda \rightarrow \infty$, $\hat{\beta}^{\text{ridge}} \rightarrow \mathbf{0}$

Bias and Variance of Ridge Regression

- When $\hat{\beta}^{\text{ols}}$ exists, we can also write

$$\begin{aligned}\hat{\beta}^{\text{ridge}} &= (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\ &= (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X}) \hat{\beta}^{\text{ols}} \\ &= \mathbf{Z} \hat{\beta}^{\text{ols}}\end{aligned}$$

where $\mathbf{Z} = (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X})$.

- How does this shrink $\hat{\beta}^{\text{ols}}$?

Bias and Variance of Ridge Regression

- The variance of $\hat{\beta}^{\text{ridge}}$ is

$$\text{Var}(\hat{\beta}^{\text{ridge}}) = (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1}$$

- The total variance $\sum_j \text{Var}(\hat{\beta}_j^{\text{ridge}})$ is a monotone decreasing function of λ .

Bias and Variance of Ridge Regression

- The the ridge estimator is biased

$$E(\hat{\beta}^{\text{ridge}}) = \mathbf{Z}\beta$$

where $\mathbf{Z} = (\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}(\mathbf{X}^T\mathbf{X})$.

- The total squared bias $\sum_j \text{Bias}^2(\hat{\beta}_j^{\text{ridge}})$ is a monotone increasing function of λ .

Understanding the Shrinkage

- Suppose we have orthogonal design matrix ($\mathbf{X}^T \mathbf{X} = n\mathbf{I}$), then $\hat{\beta}^{\text{ols}} = \frac{1}{n} \mathbf{X}^T \mathbf{y}$ and

$$\begin{aligned}\hat{\beta}^{\text{ridge}} &= (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} (\mathbf{X}^T \mathbf{X}) \hat{\beta}^{\text{ols}} \\ &= (\mathbf{I} + \lambda \mathbf{I})^{-1} \hat{\beta}^{\text{ols}} \\ &= (1 + \lambda)^{-1} \hat{\beta}^{\text{ols}},\end{aligned}$$

meaning that we just need to shrink $\hat{\beta}^{\text{ols}}$ by $(1 + \lambda)^{-1}$, i.e.,

$$\hat{\beta}_j^{\text{ridge}} = \frac{1}{1 + \lambda} \hat{\beta}_j^{\text{ols}}.$$

Understanding the Shrinkage

- $\text{Var}(\hat{\beta}_j^{\text{ridge}}) = \frac{1}{(1+\lambda)^2} \text{Var}(\hat{\beta}_j^{\text{ols}})$ (reduced from OLS!)
- $\text{Bias}(\hat{\beta}_j^{\text{ridge}}) = \frac{-\lambda}{1+\lambda} \beta_j$ (not unbiased!)
- There always exists a λ such that the MSE of $\hat{\beta}^{\text{ridge}}$ is smaller than $\hat{\beta}^{\text{ols}}$

Understanding the Shrinkage

- When the columns of \mathbf{X} are not orthogonal, let's take a singular value decomposition (SVD) of \mathbf{X} :

$$\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$$

where

- $\mathbf{U}_{n \times n}$: columns \mathbf{u}_j 's form an orthonormal basis for the column space of \mathbf{X} , $\mathbf{U}^T\mathbf{U} = \mathbf{I}$
- $\mathbf{V}_{p \times p}$: orthogonal matrix with $\mathbf{V}^T\mathbf{V} = \mathbf{I}$
- $\mathbf{D}_{n \times p}$: matrix with diagonal entries $d_1 \geq d_2 \geq \dots \geq d_p \geq 0$ being the singular values of \mathbf{X}
- Sometimes we can write $\mathbf{X} = \mathbf{F}\mathbf{V}^T$ where each columns of $\mathbf{F}_{n \times p} = \mathbf{U}\mathbf{D}$ is the so-called **principal components** and each column of \mathbf{V} is a **principal direction**.

Understanding the Shrinkage

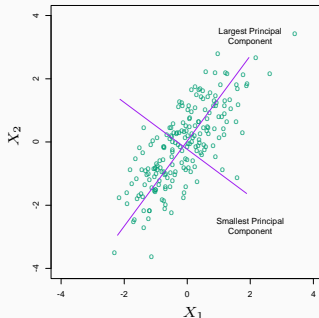


FIGURE 3.9. *Principal components of some input data points. The largest principal component is the direction that maximizes the variance of the projected data, and the smallest principal component minimizes that variance. Ridge regression projects \mathbf{y} onto these components, and then shrinks the coefficients of the low-variance components more than the high-variance components.*

Understanding the Shrinkage

- We can view PCA as (assuming \mathbf{X} centered)

$$\widehat{\Sigma} = \frac{1}{n} \mathbf{X}^T \mathbf{X} = \frac{1}{n} \mathbf{V} \mathbf{D}^2 \mathbf{V}^T$$

where $\mathbf{D}^2 = \text{diag}(d_1^2, d_2^2, \dots, d_p^2)$.

- The j th **principal component** is $\mathbf{z}_j = \mathbf{X} \mathbf{v}_j = d_j \mathbf{u}_j$ with $\text{Var}(\mathbf{z}_j) = d_j^2$.
- \mathbf{u}_j is the normalized j th principal component of \mathbf{X}
- The Ridge estimate $\widehat{\mathbf{y}}^{\text{ridge}}$ is

$$\mathbf{X} \widehat{\boldsymbol{\beta}}^{\text{ridge}} = \mathbf{X} (\mathbf{X}^T \mathbf{X} + n\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} = \sum_{j=1}^p \mathbf{u}_j \left(\frac{d_j^2}{d_j^2 + n\lambda} \mathbf{u}_j^T \mathbf{y} \right)$$

Understanding the Shrinkage

- Hence, Ridge regression can be understood as
 - (1) Perform principle component analysis of \mathbf{X}
 - (2) Project \mathbf{y} onto the principal components: $\mathbf{u}_j^T \mathbf{y}$ for each j
 - (3) Shrink the projections by the factor $d_j^2 / (d_j^2 + n\lambda)$
- Directions with smaller eigenvalues d_j^2 get more shrinkage.
- The final ridge estimate of \mathbf{y} is a sum of the p shrunk projections.

Degrees of Freedom for Ridge Regression

- Although $\hat{\beta}^{\text{ridge}}$ is p -dimensional, it does not use the full potential of the p covariates due to the shrinkage.
- For example, when $\lambda \rightarrow \infty$, all the parameter estimates are shrunk to 0. Intuitively, the d.f. is almost 0.
- If λ is 0, then it reduces to the OLS with d.f. = p
- The d.f. of a Ridge regression is between 0 and p

Degrees of Freedom for Ridge Regression

- Recall our definition of **degrees of freedom** (d.f.) in the k NN example:

$$\text{df}(\hat{f}) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Cov}(\hat{y}_i, y_i) = \frac{1}{\sigma^2} \sum_{i=1}^n \text{Trace}(\text{Cov}(\hat{\mathbf{Y}}, \mathbf{Y}))$$

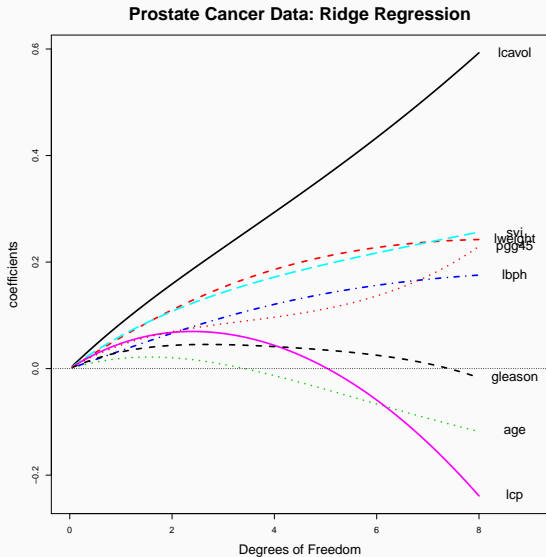
- For Ridge regression, we have

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

- Then the effective d.f. is

$$\text{df}(\lambda) = \text{Trace}(\mathbf{X}(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T) = \sum_{j=1}^p \frac{d_j^2}{d_j^2 + n\lambda}$$

Prostate Cancer Example



Selecting the Tuning Parameter λ

- The R command `lm.ridge` (from MASS package) returns GCV, which can be used to select λ .
- `glmnet` can also fit Ridge regression by setting $\alpha = 0$
- The leave-one-out cross-validation (CV) error? In the context of linear regression
 - 1 Hold the i th sample (x_i, y_i) as a test sample, fit a regression model based on the remaining $(n - 1)$ observations, and denote the coefficient as $\hat{\beta}_{[-i]}$
 - 2 Calculate the prediction error on the holdout sample $(y_i - x_i^T \hat{\beta}_{[-i]})^2$
 - 3 Repeat for every sample and

$$\text{CV} = \sum_{i=1}^n (y_i - x_i^T \hat{\beta}_{[-i]})^2$$

Selecting the Tuning Parameter λ

- In LS, we do not need to run n regression models to calculate the leave-one-out CV

$$\begin{aligned}\text{CV} &= \sum_{i=1}^n (y_i - x_i^\top \hat{\beta}_{[-i]})^2 \\ &= \sum_{i=1}^n \left(\frac{y_i - x_i^\top \hat{\beta}}{1 - \mathbf{H}_{ii}} \right)^2\end{aligned}$$

where \mathbf{H}_{ii} is the (i, i) -th entry of the projection matrix \mathbf{H} .

- Hence, we only need to run LS once and rescale the residuals.

Selecting the Tuning Parameter λ

- For Ridge regression, it is very similar

$$\text{CV}(\lambda) = \sum_{i=1}^n \left(\frac{y_i - x_i^\top \hat{\beta}_\lambda^{\text{ridge}}}{1 - \mathbf{S}_\lambda(i, i)} \right)^2$$

where $\mathbf{S}_\lambda(i, i)$ is the (i, i) -th entry of the projection matrix

$$\mathbf{S}_\lambda = \mathbf{X}(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top$$

- A modified version is called GCV (generalized CV)

$$\text{GCV}(\lambda) = \frac{\sum_{i=1}^n (y_i - x_i^\top \hat{\beta}_\lambda^{\text{ridge}})^2}{(n - \text{Trace}(\mathbf{S}_\lambda))^2}$$

- The Ridge regression solution can be viewed from a Bayesian prospective, where we give a prior distribution $\beta \sim \mathcal{N}(0, \sigma^2/\lambda)$.
- Then the posterior distribution of β is normal, with posterior mean

$$\left(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I}\right)^{-1}\mathbf{X}^T\mathbf{y},$$

and posterior variance

$$\sigma^2(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}\mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X} + \lambda\mathbf{I})^{-1}$$

The solution is not invariant with respect to the scale of the predictors!

we normalize the columns of the design matrix \mathbf{X} such that they have unit sample variance. We further center the data, that is, both y and the columns of \mathbf{X} have mean zero. Then, we can fit a linear regression model without an intercept (we don't penalize the intercept). The parameters on the original scale can be reversely solved.

Some packages (e.g. “`glmnet`”) in `R` handles the centering and scaling automatically: it will do the transformation before running the algorithm, and then will transform the obtained results back to the original scale.

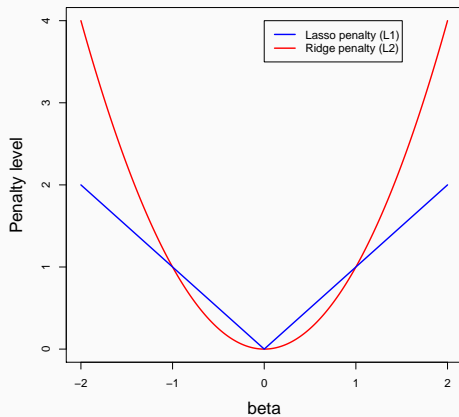
Lasso: Least Absolute Shrinkage and Selection Operator

- The Ridge regression shrinks the coefficients towards 0, however, they are not exactly zero. Hence, we haven't achieve any "selection" of variables.
- Parsimony: we would like to select a small subset of predictions. Forward/backward/subset does not provide global solution and can be myopic at each step.
- Lasso provides a continuous process. We will discuss:
 - The formulation, the solution when \mathbf{X} is orthogonal
 - Computation methods and solution path

Least absolute shrinkage and selection operator (Tibshirani 1996)

$$\arg \min_{\beta} \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \|\beta\|_1$$

- Shrinkage of the ℓ_1 norm of the parameters
- Selection of parameters, some will be exactly 0





Lasso Under Orthogonal Design

Again, it will be helpful to view Lasso assuming orthogonal design, i.e., $\mathbf{X}^T \mathbf{X} = n\mathbf{I}_p$. Then

$$\begin{aligned}\frac{1}{n} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 &= \frac{1}{n} \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} + \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{X}\boldsymbol{\beta}\|^2 \\ &= \frac{1}{n} \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}}\|^2 + \frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{X}\boldsymbol{\beta}\|^2\end{aligned}$$

where the cross product term

$$2(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}})^T (\mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{X}\boldsymbol{\beta}) = 2\mathbf{r}^T (\mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{X}\boldsymbol{\beta}) = 0,$$

since the second term is in the column space of \mathbf{X} , while \mathbf{r} is orthogonal to that space.

Lasso Under Orthogonal Design

- Since $\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}}\|^2$ is not a function of $\boldsymbol{\beta}$, we minimize

$$\frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1$$

- Then, we have

$$\begin{aligned}\hat{\boldsymbol{\beta}}^{\text{lasso}} &= \arg \min_{\boldsymbol{\beta}} \frac{1}{n} \|\mathbf{X}\hat{\boldsymbol{\beta}}^{\text{ols}} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \arg \min_{\boldsymbol{\beta}} \frac{1}{n} (\hat{\boldsymbol{\beta}}^{\text{ols}} - \boldsymbol{\beta})^\top \mathbf{X}^\top \mathbf{X} (\hat{\boldsymbol{\beta}}^{\text{ols}} - \boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \arg \min_{\boldsymbol{\beta}} (\hat{\boldsymbol{\beta}}^{\text{ols}} - \boldsymbol{\beta})^\top (\hat{\boldsymbol{\beta}}^{\text{ols}} - \boldsymbol{\beta}) + \lambda \|\boldsymbol{\beta}\|_1 \\ &= \arg \min_{\boldsymbol{\beta}} \sum_{j=1}^p (\hat{\beta}_j^{\text{ols}} - \beta_j)^2 + \lambda |\beta_j|.\end{aligned}$$

- This means we can solve the lasso estimators individually from the OLS estimator.

Lasso Under Orthogonal Design

- Each of the β_j 's is essentially solving for

$$\arg \min_x (x - a)^2 + \lambda|x|, \quad \lambda > 0$$

- The solution is simply

$$\begin{aligned} \hat{\beta}_j^{\text{lasso}} &= \begin{cases} \hat{\beta}_j^{\text{ols}} - \lambda/2 & \text{if } \hat{\beta}_j^{\text{ols}} > \lambda/2 \\ 0 & \text{if } |\hat{\beta}_j^{\text{ols}}| \leq \lambda/2 \\ \hat{\beta}_j^{\text{ols}} + \lambda/2 & \text{if } \hat{\beta}_j^{\text{ols}} < -\lambda/2 \end{cases} \\ &= \text{sign}(\hat{\beta}_j^{\text{ols}}) \left(|\hat{\beta}_j^{\text{ols}}| - \lambda/2 \right)_+ \\ &\doteq \text{SoftTH}(\beta_j^{\text{ols}}, \lambda) \end{aligned}$$

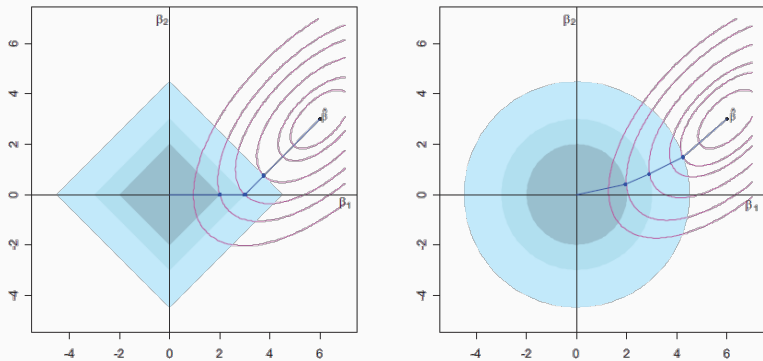
- A large λ will shrink some of the coefficients to exactly zero, which achieves “variable selection”.

- The Lasso optimization problem is equivalent to

$$\begin{aligned} & \underset{\beta}{\text{minimize}} && \sum_{i=1}^n \left(y_i - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \\ & \text{subject to} && \sum_{j=1}^p |\beta_j| \leq s \end{aligned}$$

- Each value of λ corresponds to a unique value of s .
- Compare Ridge and Lasso?

Linear Regression



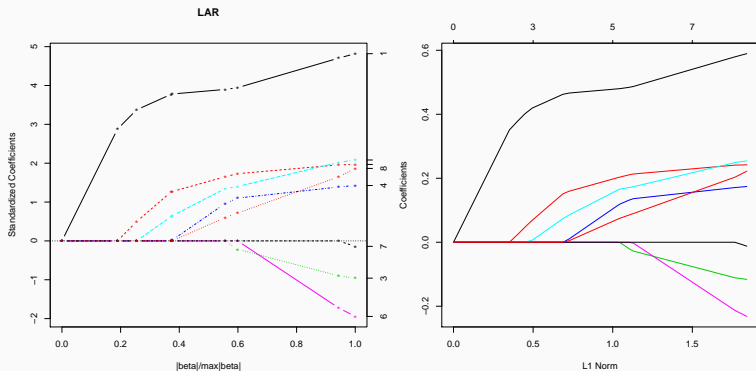
Comparing Lasso and Ridge solutions

Computation of Lasso Solution

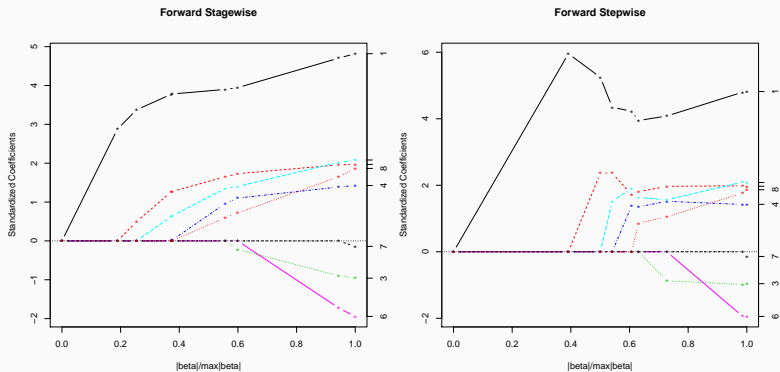
- The Lasso problem is convex, although it may not be strictly convex in β when p is large
- The solution is a global minimum, but may not be the **unique** global one
- The Lasso solution is **unique under conditions of the covariance matrix**

Computation of Lasso Solution

- Shooting algorithm (Fu 1998): sequentially and iteratively update each parameter estimate (coordinate descent algorithm).
- Least angle regression (Efron et al. 2004)
 - The path of solutions is piecewise linear in λ
 - Cost is approximately one least-squares calculation $\mathcal{O}(np^2)$
 - Connection with stagewise regression
- Coordinate descent (Friedman et al 2010): The most popular implementation, `glmnet` package; $\mathcal{O}(np)$
 - Also provides the solution path for the entire sequence of λ , starting with the largest one
 - Use the previous estimation of β as a warm start for smaller λ



Comparing least angle regression with coordinate descent



Comparing stagewise regression with stepwise regression

- Ridge is ℓ_2 penalty
- Lasso is ℓ_1 penalty
- Best subset is ℓ_0 penalty
- Bridge penalty is ℓ_q normal

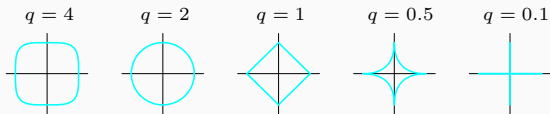


FIGURE 3.12. *Contours of constant value of $\sum_j |\beta_j|^q$ for given values of q .*

- Elastic-net is a hybrid of ℓ_1 and ℓ_2 :

$$\lambda_1 \|\boldsymbol{\beta}\|_1 + \lambda_2 \|\boldsymbol{\beta}\|_2^2$$

R Functions

- Use `R` help and `R` manuals
- Linear models: function `lm`
- QR decomposition `qr`; Cholesky decomposition `chol`; PCA `princomp`, `prcomp`; SVD `svd`.
- Ridge regression:
 - package `MASS`; function `lm.ridge`
 - package `glmnet`; function `glmnet` and `cv.glmnet` with `alpha = 0`
- Lasso:
 - package `lars`; function `lars`
 - package `glmnet`; function `glmnet` and `cv.glmnet` with `alpha = 1`